

---

# Applying Time-Series Causal Discovery to Understand Algal Bloom Mechanisms

---

Ayush Prasad  
NatureDots

Snehal Verma  
NatureDots

Mohammad Aatish Khan  
NatureDots

firstname.lastname@naturedots.com

## Abstract

1 Harmful algal blooms present a significant ecological and public health challenge,  
2 yet their drivers are embedded in complex, time-lagged systems with unobserved  
3 confounding variables. Traditional correlational methods often fail to uncover the  
4 true causal structure from observational time-series data. In this work, we apply  
5 a suite of modern constraint-based causal discovery algorithms to a real-world  
6 dataset of daily water quality and weather from Lake Mendota. By explicitly  
7 modeling time lags and allowing for the presence of latent confounders using  
8 the Time-Series Iterative Causal Discovery (TS-ICD) algorithm, our approach  
9 uncovers a scientifically plausible causal graph. We demonstrate the necessity of  
10 this approach by comparing its structural output and computational complexity  
11 against a baseline algorithm (PC) that assumes causal sufficiency. Our results  
12 identify key drivers of chlorophyll-a and reveal evidence of confounding between  
13 same-day water temperature and algae levels.

## 14 1 Introduction

15 Harmful algal blooms (HABs) are a growing global concern, posing threats to aquatic ecosystems,  
16 local economies, and public health. Studies show that the frequency of HABs is increasing in diverse  
17 environments, from high-latitude coastal waters affected by warming and freshening (8) to inland  
18 temperate lakes like our study site, Lake Mendota, which has a long history of eutrophication driven  
19 by nutrient runoff (1). Understanding the causal mechanisms that trigger these blooms is crucial.  
20 While factors like nutrient loading, temperature, and sunlight are known drivers, their interactions are  
21 complex, involving significant time delays and confounding from unmeasured variables.

22 Analyzing observational time-series data from such systems is a significant challenge. While modern  
23 predictive approaches, such as physics-guided machine learning, have shown success in forecasting  
24 water quality dynamics (5), they are primarily focused on prediction rather than discovery of the  
25 underlying causal structure. Early causal approaches in Earth sciences often relied on methods like  
26 Granger causality or Dynamic Bayesian Networks (DBNs) (4), which can fail to distinguish direct  
27 from indirect causes and are highly susceptible to being misled by latent confounders (2).

28 To address these limitations, we apply modern causal discovery techniques to the Lake Mendota  
29 dataset (10). Our primary method is the Time-Series Iterative Causal Discovery (TS-ICD) algorithm  
30 (6), which is specifically designed for time-series data and is robust to the presence of latent con-  
31 founders. By comparing its output to simpler algorithms like PC (9), we demonstrate the critical  
32 importance of selecting a method whose assumptions align with the realities of the system under  
33 study.

34 Our contributions are:

1. We present an application of a recently developed suite of causal discovery algorithms to a real-world limnological dataset, overcoming limitations of prior causal approaches.
2. We conduct a comparative analysis showing that explicitly modeling latent variables and leveraging temporal structure yields more efficient and plausible causal models.
3. The resulting causal graph provides new, testable hypotheses about the specific time-lagged drivers of algal blooms in Lake Mendota.

## 2 Methodology

We developed our methodology to discover and compare causal structures from observational time-series data. We first describe our dataset and pre-processing steps, then provide a detailed overview of the three constraint-based causal discovery algorithms used in our comparative analysis, also listing their different underlying assumptions.

### 2.1 Data and Pre-processing

We use a high-frequency, daily resolution dataset of water quality and meteorology for Lake Mendota, USA, sourced from the LakeBeD-US benchmark dataset (10). For our analysis, we selected a focused set of five variables known to be relevant to algal blooms: `precipitation_sum`, `temperature_2m_mean`, `shortwave_radiation_sum`, `water temp`, and `chl_a_rfu` (chlorophyll-a). To account for seasonality as a potential confounder, we engineered cyclical features from the date, `doy_sin` and `doy_cos`. To model time-lagged relationships, we transformed the data into a static format using a lag window of 4 days, meaning an outcome at time  $t$  is predicted by variables from  $t - 1$  through  $t - 4$ . All conditional independence (CI) queries were performed using a partial correlation test (`CondIndepParCorr`) with a significance level of  $\alpha = 0.01$ .

### 2.2 Causal Discovery Algorithms

The PC (Peter-Clark) algorithm serves as our baseline, operating under the strong assumption of *causal sufficiency* (no latent confounders) to recover the causal skeleton and equivalence class (CPDAG) of the underlying graph (9). It begins with a fully connected undirected graph and first performs an adjacency search, iteratively removing edges between variables  $X$  and  $Y$  if a separating set  $S$  is found such that  $(X \perp Y | S)$ . After establishing the skeleton, it orients edges by first identifying v-structures ( $X \rightarrow Z \leftarrow Y$ ) where the separating set  $S_{XY}$  does not contain the collider  $Z$ , and then exhaustively applying a set of orientation rules.

Our primary method is the Time-Series Iterative Causal Discovery (TS-ICD) algorithm, designed for time-series and *not assuming causal sufficiency*. TS-ICD outputs a Partial Ancestral Graph (PAG), which can represent confounding via bidirected edges ( $\leftrightarrow$ ). The algorithm initializes a fully connected PAG with uncertain endpoints ( $o - o$ ) and proceeds iteratively, performing adjacency removal and orientation phases. Critically, TS-ICD leverages the temporal structure by enforcing the arrow of time and assuming stationarity, significantly pruning the search space.

We also evaluate the Ordered Iterative Causal Discovery (OrdICD) algorithm to demonstrate the benefit of incorporating background knowledge. OrdICD also allows for latent confounders but gains efficiency by using a known causal order ( $\prec$ ). Given the natural temporal ordering in our unrolled data, where variables at  $t - k$  precede those at  $t - j$  for  $k > j$ , the algorithm constrains its search for a separating set, which dramatically reduces the number of required CI tests.

## 3 Results and Analysis

Our experimental analysis is presented in two parts. First, we evaluate the computational performance of the three causal discovery algorithms. Second, we present and interpret the learned causal graph for Lake Mendota from our primary method, TS-ICD, and contrast its structure with the expected output from the baseline.

### 3.1 Computational Performance

We measured the computational complexity of each algorithm by the total number of CI tests required. The results (Table 1) show a clear performance hierarchy. The baseline PC algorithm was the most exhaustive (2625 tests), while OrdICD (1091 tests) and especially TS-ICD (634 tests) were significantly more efficient.

Table 1: Computational cost of causal discovery methods measured by the number of conditional independence (CI) tests.  $k$  is the conditioning order (size of the conditioning set). PC performs the most tests. ICD-Ord reduces tests by using background temporal order: TS-ICD requires the fewest especially at higher orders ( $k \geq 2$ ).

Algorithm	Total CI Tests	k=0	k=1	k=2
PC	2625	276	1168	1181
OrdICD	1091	276	671	144
TS-ICD	634	276	319	39

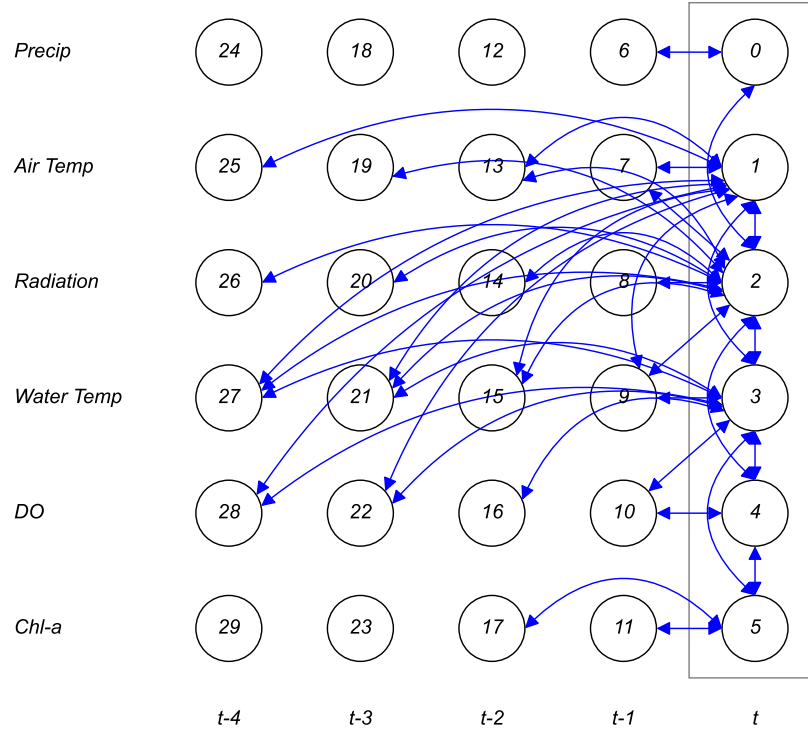


Figure 1: Partial Ancestral Graph (PAG) learned by TS-ICD for Lake Mendota. Columns show lags  $t-4$  to  $t$ ; rows list variables (Precip, Air Temp, Radiation, Water Temp, DO, Chl-a). Solid arrows ( $\rightarrow$ ) denote putative direct causal effects; bidirected edges ( $\leftrightarrow$ ) at  $t$  indicate possible latent confounding. From the graph we infer: (i) state persistence ( $\text{Chl-a}_{t-1} \rightarrow \text{Chl-a}_t$  and  $\text{Water Temp}_{t-1} \rightarrow \text{Water Temp}_t$ ); (ii) a lagged light effect  $\text{Radiation}_{t-1} \rightarrow \text{Chl-a}_t$  consistent with diel growth; (iii) a meteorological-physical link  $\text{Air Temp}_{t-1} \rightarrow \text{Water Temp}_t$ ; and (iv) a contemporaneous  $\text{Water Temp}_t \leftrightarrow \text{Chl-a}_t$  dependence, suggesting an unobserved regime (e.g., calm, sunny conditions) influencing both.

### 85 3.2 Causal Structure of Algal Bloom Drivers

86 The PAG learned by the TS-ICD algorithm for Lake Mendota (Figure 1) summarizes the inferred  
87 relationships among drivers and chlorophyll-a. This model provides a data-driven hypothesis of the  
88 dynamics driving chlorophyll-a and is consistent with basic lake physics and phytoplankton ecology.  
89 First, the graph recovers state persistence (e.g.,  $\text{Chl-a}_{t-1} \rightarrow \text{Chl-a}_t$  and  $\text{temp}_{t-1} \rightarrow \text{temp}_t$ ), which  
90 is expected because biomass accumulates over multiple days and the surface layer has thermal inertia.  
91 Second, we observe a direct, time-lagged link from shortwave radiation at  $t-1$  to chlorophyll-a at  
92  $t$ . At a daily resolution this lag is plausible: light drives photosynthesis within hours, but changes  
93 in measured chlorophyll (via pigments and cell division) typically appear after one diel cycle. The  
94 graph also shows the meteorological-physical pathway in which air temperature at  $t-1$  points to  
95 water temperature at  $t$ , reflecting surface heat fluxes that warm and stabilize the mixed layer.

96 The remaining meteorological drivers act in ways that agree with known mechanisms. Edges from  
97 wind speed point into physical variables and, where identified, into chlorophyll. Stronger winds  
98 deepen the mixed layer, dilute surface biomass, and can lower near-surface fluorescence on short  
99 timescales; they can also resuspend nutrients that influence growth with a delay. Precipitation links  
100 appear as upstream inputs consistent with runoff and nutrient pulses. Where orientation is weak  
101 at daily resolution, the PAG leaves an open endpoint, which correctly records that the data do not  
102 compel a single direction.

103 Most importantly, the PAG contains a bidirected edge between same-day water temperature and  
104 chlorophyll-a ( $\text{temp}_t \leftrightarrow \text{Chl-a}_t$ ), indicating a latent common cause. This is a plausible representation  
105 of weather regimes that are only partially measured in our covariates. Calm, sunny conditions can  
106 simultaneously warm the surface layer and increase light exposure by reducing mixing, producing a  
107 same-day association between temperature and chlorophyll that is not purely causal in either direction.  
108 By allowing a bidirected mark, TS-ICD attributes this dependence to an unobserved factor rather than  
109 forcing a direct arrow.

110 This interpretation also explains the difference to the PC baseline. Because PC assumes causal  
111 sufficiency, it must explain the same-day temperature-chlorophyll correlation with a directed edge,  
112 which is likely spurious in the presence of unmeasured regime variables. TS-ICD respects temporal  
113 ordering, uses only admissible past variables in separating sets, and reports uncertainty through PAG  
114 edge marks.

## 115 4 Conclusion

116 In this work, we presented an exploratory study on the application of modern time-series causal  
117 discovery algorithms to a real-world observational dataset. Our comparative analysis showed that  
118 the specialized TS-ICD algorithm, which accounts for latent confounders and temporal structure,  
119 was significantly more computationally efficient than baseline methods. The resulting causal model  
120 identified plausible time-lagged drivers and, crucially, highlighted a confounded relationship between  
121 two contemporaneous variables that was missed by the baseline PC algorithm. The linear assumption  
122 of the partial correlation test is a simplification; future work could explore the use of non-parametric  
123 or kernel-based conditional independence tests to capture potentially nonlinear relationships. Our  
124 analysis also relied on an assumption of stationarity. A more advanced approach would be to  
125 investigate methods that can accommodate non-stationary or regime-shifting dynamics in the data.  
126 Furthermore, the output of our analysis is a qualitative causal graph. A valuable next step could also  
127 be to use the learned Partial Ancestral Graph (PAG) as a structural prior to fit a parametric model,  
128 such as a Structural Vector Autoregression, to quantify the strengths of the discovered causal effects.  
129 This would also provide the foundation for performing counterfactual analysis.

## 130 References

- 131 [1] S. R. Carpenter, J. J. Cole, M. L. Pace, et al. Trophic cascades, nutrient trading, and lake  
132 eutrophication. *Ecosystems*, 11(5):781–796, 2008.
- 133 [2] Michael Eichler. Causal inference in time series analysis. In *Handbook of Causal Analysis for*  
134 *Social Research*, pages 327–354. Springer, 2013.

- 135 [3] A. Gerhardus and J. Runge. High-recall causal discovery for autocorrelated time series with  
136 latent confounders. In *Advances in Neural Information Processing Systems*, 33, 2020.
- 137 [4] Jaroslav Hlinka, D. Hartman, M. Vejmelka, et al. Reliability of causal discovery from subsampled  
138 time series. In *Chaos: An Interdisciplinary Journal of Nonlinear Science*, 27(12):126708, 2017.
- 139 [5] X. Jia, J. S. Read, J. P. Zwart, et al. Physics-guided machine learning for predicting lake water  
140 quality. In *Communications Earth & Environment*, 6(1):1–11, 2025.
- 141 [6] Raanan Yehezkel Rohekar, Shami Nisimov, Yaniv Gurwicz, and Gal Novik. From temporal to  
142 contemporaneous iterative causal discovery in the presence of latent confounders. In *International  
143 Conference on Machine Learning*, pages 39939–39950. PMLR, 2023.
- 144 [7] Raanan Rohekar, Shami Nisimov, Yaniv Gurwicz, and Gal Novik. Iterative causal discovery in  
145 the possible presence of latent confounders and selection bias. In *Advances in Neural Information  
146 Processing Systems*, 34, 2021.
- 147 [8] Edson Silva, François Counillon, Julien Brajard, et al. Warming and freshening coastal waters  
148 impact harmful algal bloom frequency in high latitudes. In *Communications Earth & Environment*,  
149 6(1):445, 2025.
- 150 [9] Peter Spirtes, Clark N. Glymour, and Richard Scheines. *Causation, Prediction, and Search*. MIT  
151 press, 2000.
- 152 [10] Y. Yabansu and P. G. Hanson. LakeBeD-US: a benchmark dataset of lake water quality and  
153 meteorology for 32 US lakes. In *Earth System Science Data*, 17, 3141–3157, 2025.